



**Europäisches
Patentamt**

**European
Patent Office**

**Office européen
des brevets**

REC'D 13 DEC 2004

WIPO

PCT

Bescheinigung

Certificate

Attestation

Die angehefteten Unterla-
gen stimmen mit der
ursprünglich eingereichten
Fassung der auf dem näch-
sten Blatt bezeichneten
europäischen Patentanmel-
dung überein.

The attached documents
are exact copies of the
European patent application
described on the following
page, as originally filed.

Les documents fixés à
cette attestation sont
conformes à la version
initialement déposée de
la demande de brevet
européen spécifiée à la
page suivante.

Patentanmeldung Nr. Patent application No. Demande de brevet n°

03090256.3

**PRIORITY
DOCUMENT**

SUBMITTED OR TRANSMITTED IN
COMPLIANCE WITH RULE 17.1(a) OR (b)

Der Präsident des Europäischen Patentamts;
Im Auftrag

For the President of the European Patent Office

Le Président de l'Office européen des brevets
p.o.

R C van Dijk

THIS PAGE BLANK (USPTO)



Anmeldung Nr:
Application no.: 03090256.3
Demande no:

Anmeldetag:
Date of filing: 19.08.03
Date de dépôt:

Anmelder/Applicant(s)/Demandeur(s):

FRAUNHOFER-GESELLSCHAFT ZUR FÖRDERUNG DER
ANGEWANDTEN FORSCHUNG E.V.
Leonrodstrasse 68
80636 München
ALLEMAGNE

Bezeichnung der Erfindung/Title of the invention/Titre de l'invention:
(Falls die Bezeichnung der Erfindung nicht angegeben ist, siehe Beschreibung.
If no title is shown please refer to the description.
Si aucun titre n'est indiqué se référer à la description.)

A method for online detection and classification of anomalous objects

In Anspruch genommene Priorität(en) / Priority(ies) claimed / Priorité(s)
revendiquée(s)
Staat/Tag/Aktenzeichen/State/Date/File no./Pays/Date/Numéro de dépôt:

Internationale Patentklassifikation/International Patent Classification/
Classification internationale des brevets:

H04L12/26

Am Anmeldetag benannte Vertragsstaaten/Contracting states designated at date of
filing/Etats contractants désignées lors du dépôt:

AT BE BG CH CY CZ DE DK EE ES FI FR GB GR HU IE IT LU MC NL
PT RO SE SI SK TR LI

THIS PAGE BLANK (USPTO)

EPO-BERLIN
19-08-2003

System and method for online detection and classification of anomalous objects

1 Field of the invention

This invention relates generally to detection and classification of anomalous objects in continuous data streams, and, in particular, to detection and classification of anomalous packets in computer networks, of anomalous records in computer logs, of anomalous measurements of other activities such as operation of mechanical devices, chemical, biological and medical processes.

2 Summary of the invention

A system and method are disclosed for online detection and classification of anomalous objects in continuous data streams. The system consists of the online anomaly detection engine and of the optional feature extraction and classification components.

The online anomaly detection engine consists of a processing unit having memory for storing the incoming data, the limited working set, and the geometric representation of the normal (non-anomalous) data objects by means of a parametric hypersurface; stored programs including the programs for processing of incoming data; and a processor controlled by the stored programs. The processor includes the components for construction and update of the geometric representation of normal data objects, and for the detection of anomalous objects based on the stored representation of normal data objects.

The component for construction and update of the geometric representation receives a data entry and imports it into the representation such that the smallest volume enclosed by the hypersurface and consistent with the pre-defined expected fraction of anomalous objects is maintained; the component further identifies the least relevant entry in the working set and removes it while maintaining the smallest volume enclosed by the hypersurface. Detection of the anomalous objects is performed by checking if the objects fall within or outside of the hypersurface representing the normality.

As an embodiment of the invention, the architecture of the system for detection and classification of computer intrusions is disclosed. The system consists of the feature extraction component receiving data from the audit stream; of the online anomaly detection engine; and of the classification component, produced by the event learning engine trained on the database of appropriate events.

The invention is also applicable to monitoring of the measurements of physical parameters of operating mechanical devices, of the measurements of chemical processes and of the measurement of biological activity.

EPO-BERLIN

19-08-2003

3 Brief description of the drawings

~~Fig. 1~~ is the schematic diagram of the system.

Fig. 2 is the flow-chart of the algorithm utilized for the construction and the update of the geometric representation of normal objects.

Fig. 3 is the schematic diagram of the system for detection and classification of computer intrusions based on the disclosed online anomaly detection engine.

4 Claims

1. A method for online detection of anomalous objects in a continuous data stream based on the following principles:
 - a) geometric representation of the class of the normal objects,
 - b) dynamic adaptation of the said representation upon the arrival on the new data.
2. The method as defined in claim 1, wherein the objects are data packets in communication networks or representations thereof.
3. The method as defined in claim 1, wherein the objects are the records obtained from logging the processes on the computer system or representations thereof.
4. The method as defined in claim 1, wherein the objects are measurements of physical characteristics of operating mechanical devices.
5. The method as defined in claim 1, wherein the objects are measurements of biological activity.
6. The method as defined in claims 1-5, wherein the class of the normal objects is represented by a parametric boundary hypersurface in a pre-defined feature space, such that said surface encloses the smallest volume among all possible surfaces consistent with the pre-defined fraction of the anomalous objects.
7. The method as defined in claim 6, wherein the feature space is induced by a suitably-defined similarity function between the data objects satisfying the conditions under which the said function acts as an inner product in the said feature space.
8. The method as defined in claim 7, wherein the anomalous objects are determined as the ones lying on the other side of the boundary hypersurface representing the class of the normal objects.

9. The method as defined in claim 8, wherein the representation of the class of the normal objects is updated upon the arrival of a new data object, normal or anomalous.
10. The method as defined in claim 9, wherein the update of the representation of the class of the normal objects comprises the adjustment of parameters of said representation so as to incorporate the new object while maintaining the optimality of the representation.
11. The method as defined in claim 9, wherein the update of the representation of the class of the normal objects comprises the adjustment of parameters of said representation so as to remove the least-relevant object while maintaining the optimality of the representation.
12. The method to initialize the method defined in claim 9, so that the smallest-volume representation is maintained starting from the moment when the smallest possible amount of data is available for which such representation is possible given the pre-defined fraction of anomalous data.
13. A method for detection of anomalous objects in a collected database whereby the method as defined in claim 9 is applied to the database one entry at a time.
14. The method as defined in claim 13, wherein the objects in the database are collected data packets in communication networks or representations thereof.
15. The method as defined in claim 13, wherein the objects in the database are the collected records obtained from logging the processes on the computer system or representations thereof.
16. The method as defined in claim 13, wherein the objects in the database are collected measurements of physical characteristics of operating mechanical devices.
17. The method as defined in claim 13, wherein the objects in the database are collected measurements of biological activity.
18. A method for classification of anomalous objects whereby the methods as defined in claims 9 or 13 are used as a pre-processing step.
19. A method for classification of anomalous objects whereby the information collected by the detection methods as defined in claims 9 or 13 is combined with the information used by other classification methods.

EPO-BERLIN

19-08-2003

5 Detailed description of the invention

The invention constitutes the system and method for online detection and classification of anomalous objects. The main component of the invention is the online anomaly detection engine.

The overall scheme of the said system is depicted in Fig. 1. The input of the system is a data stream D1 pertaining to a particular application. The data stream can be packets in communication networks, entries in the various activity logs in computer systems, measurements of physical characteristics of operating mechanical devices, measurements of parameters of chemical processes, measurements of biological activity, and others. The crucial feature of the invention is that it can deal with *continuous* data streams in an *online* fashion. If for some reason online data processing is not desired, or if some intermediate storage of the data in a database or in a storage buffer is required, the invention can likewise be applied to the stored data, as indicated by D2, by sequentially processing stored entries. Each of the incoming data entries is supplied to the feature extraction unit A1, which performs the pre-processing required to obtain the features D3 relevant for a particular application. Alternatively, if the data entries are such that they can be directly used in a detection/classification method, the feature extraction step can be skipped.

The main step A2 of the online anomaly detection engine consists of the construction and of the update of the geometric representation of the notion of normality. The geometric representation of normality D4 is a parametric hypersurface enclosing the smallest volume among all possible surfaces consistent with the pre-defined fraction of the anomalous objects. The said hypersurface is constructed in the feature space induced by a suitably-defined similarity function between the data objects ("kernel function") satisfying the conditions under which the said function acts as an inner product in the said feature space ("Mercer conditions"). The update of the said representation involves the adjustment so as to incorporate the latest data entry, and the adjustment so as to remove the least relevant data entry so as to retain the encapsulation of the smallest volume. Once the normality representation is updated, the anomaly detection A3 is performed by assigning to the data entry the status of a normal entry, if the entry falls into the volume encompassed by the normality representation, or the status of an anomalous entry, if the entry lies outside of the volume encompassed by the normality representation.

The output of the online anomaly detection engine is used to issue the anomaly warning D5 and/or to trigger the classification component A4 which can utilize any known classification method such as decision trees, neural networks, support vector machines, Fischer discriminant etc. The geometric representation of normality can also be supplied to the classification component if this is required by the method.

In the exemplary embodiment of the step A2 the hypersurface representing the class of normal events is represented by the set of parameters x_1, \dots, x_n , one for each entry

+49 30 8825823

in the working set. The size n of the working set is chosen in advance by the user. The parameters are further restricted to be non-negative, and to have values less than or equal to $C = 1/(n\nu)$, where ν is the expected fraction of the anomalous events in the data stream, to be set by the user. The working set is partitioned into the "set O " of the entries whose parameters x_i are equal to zero, "set E " of the entries whose parameters x_i are equal to C , and the "set S " of the remaining entries. The operation of step A2 is illustrated in Fig. 2. Upon the arrival of the data entry k , the following three main actions are performed: in step A2.5 the data entry is "imported" into the working set, in step A2.6 the least relevant data entry l is sought in the working set, and in step A2.7 the data entry l is removed from the working set. The importation and removal operations maintain the minimal volume enclosed by the hypersurface and consistent to the pre-defined expected fraction of anomalous objects. These operations are explained in more detail below. The relevance of the data entry can be judged either by the time stamp on the entry or by the value of parameter x_i assigned to the entry. The steps A2.1–A2.4 are the initialization operations to be performed when not enough data entries have been observed in order to bring the system into equilibrium.

Construction of the hypersurface $D4$ enclosing the smallest volume and consistent with the pre-defined expected fraction of anomalous objects amounts, as shown in the article "Support Vector Data Description" by D.M.J. Tax and R.P.W. Duin, *Pattern Recognition Letters*, vol. 20, pages 1191–1199, (1999), to solving the following mathematical programming problem:

$$\max_{\mu} \min_{\substack{0 \leq x \leq C \\ a^T x + b = 0}} : W = -c^T x + \frac{1}{2} x^T K x + \mu(a^T x + b), \quad (1)$$

where K is a $n \times n$ matrix that consists of evaluations of the given kernel function for all data points in the working set: $K_{ij} = \text{kernel}(p_i, p_j)$, c is the vector of the numbers at the main diagonal of K , a is the vector of n ones, and $b = -1$. The parameter C is related to the expected fraction of the anomalous objects. The necessary and sufficient condition for the optimality of the representation attained by the solution to problem (1) is given by the well-known Karush-Kuhn-Tucker conditions. When all the points in the working set satisfy the said conditions, the working set is said to be in equilibrium. Importation of a new data entry into, or removal of an existing data entry from a working set may result in the violation of the said conditions. In such case, adjustments of the parameters x_1, \dots, x_n are necessary, in order to bring the working set back into the equilibrium. An algorithm for performing such adjustments, based on the Karush-Kuhn-Tucker conditions, for a different mathematical programming problem — Support Vector Learning — was presented in the article "Incremental and Decremental Support Vector Learning" by G. Cauwenberghs and T. Poggio, *Advances in Neural Information Processing Systems 13*, pages 409–415, (2001). By deriving the Karush-Kuhn-Tucker conditions for the problem (1), the necessary ingredients for the

application of the method of Cauwenberghs and Poggio can be obtained.

Special care needs to be taken at the initial phase of the operation of the online anomaly detection engine as described in Fig. 2. When the number of data entries in the working set is less than or equal to $\lfloor \frac{1}{\epsilon} \rfloor$ (the greatest integer smaller than or equal to $\frac{1}{\epsilon}$), equilibrium cannot be reached and the method of Cauwenberghs and Poggio cannot be applied. The initialization steps A2.1-A2.4 are designed to handle this special case and to bring the working set into the equilibrium after the smallest possible number of data entries has been seen.

The exemplary embodiment of the online anomaly detection method in the system for detection and classification of computer intrusions is depicted in Fig. 3. The input of the said system is an audit stream which contains network packets and records in the audit logs of computers. The audit stream is input into the feature extraction component comprising a set of filters to extract the relevant features. The extracted features are read by the online anomaly detection engine which identifies anomalous objects (packets or log entries) and issues an event warning if the event is discovered to be anomalous. Classification of the detected anomalous events is performed by the classification component previously trained to classify the anomalous events collected and stored in the event database.

APPENDIX

ONLINE SVM LEARNING: FROM CLASSIFICATION TO DATA DESCRIPTION AND BACK

Abstract. The paper presents two useful extensions of the incremental SVM in the context of online learning. An online support vector data description algorithm enables application of the online paradigm to unsupervised learning. Furthermore, online learning can be used in the large-scale classification problems to limit the memory requirements for storage of the kernel matrix. The proposed algorithms are evaluated on the task of online monitoring of EEG data, and on the classification task of learning the USPS dataset with a-priori chosen working set size.

INTRODUCTION

Many real-life machine learning problems can be more naturally viewed as online rather than batch learning problems. Indeed, the data is often collected continuously in time, and, more importantly, the concepts to be learned may also evolve in time. Significant effort has been spent in the recent years on development of online SVM learning algorithms (e.g. [17, 13, 7, 12]). The elegant solution to online SVM learning is the incremental SVM [4] which provides a framework for exact online learning. In the wake of this work two extensions to the regression SVM have been independently proposed [10, 9].

One should note, however, a significant restriction on the applicability of the above-mentioned supervised online learning algorithms: the labels may not be available online, as it would require manual intervention at every update step. A more realistic scenario is the update of the existing classifier when a new batch of data becomes available. The true potential of online learning can only be realized in the context of unsupervised learning.

An important and relevant unsupervised learning problem is one-class classification [11, 14]. This problem amounts to constructing a multi-dimensional data description, and its main application is novelty (outlier) detection. In this case online algorithms are essential, for the same reasons that made online learning attractive in the supervised case: the dynamic nature of data

and drifting concepts. An online support vector data description (SVDD) algorithm based on the incremental SVM is proposed in this paper.

Looking back at the supervised learning, a different role can be seen for online algorithms. Online learning can be used to overcome memory limitations typical for kernel methods on large-scale problems. It has been long known that storage of the full kernel matrix, or even the part of it corresponding to support vectors, can well exceed the available memory. To overcome this problem, several subsampling techniques have been proposed [16, 1]. Online learning can provide a simple solution to the subsampling problem: make a sweep through the data with a limited working set, each time adding a new example and removing the least relevant one. Although this procedure results in an approximate solution, an experiment on the USPS data presented in this paper shows that significant reduction of memory requirements can be achieved without major decrease in classification accuracy.

To present the above-mentioned extensions we first need an abstract formulation of the SVM optimization problem and a brief overview of the incremental SVM. Then the details of our algorithms are presented, followed by their evaluation on real-life problems.

PROBLEM DEFINITION

A smooth extension of the incremental SVM to the SVDD can be carried out by using the following abstract form of the SVM optimization problem:

$$\max_{\mu} \min_{\substack{0 \leq \alpha \leq C \\ \alpha^T x + 1 = 0}} : W = -c^T x + \frac{1}{2} x^T K x + \mu (a^T x + b), \quad (1)$$

where c and a are $n \times 1$ vectors, K is a $n \times n$ matrix and b is a scalar. By defining the meaning of the abstract parameters a , b and c for the particular SVM problem at hand, one can use the same algorithmic structure for different SVM algorithms. In particular, for the standard support vector classifiers [19], take $c = 1$, $a = y$, $b = 0$ and the given regularization constant C ; the same definition applies to the ν -SVC [15] except that $C = \frac{1}{\nu\gamma}$; for the SVDD [14, 18], the parameters are defined as: $c = \text{diag}(K)$, $a = y$ and $b = -1$.

Incremental (decremental) SVM provides a procedure for adding (removing) one example to (from) an existing optimal solution. When a new point k is added, its weight α_k is initially assigned to 0. Then the weights of other points and μ should be updated, in order to obtain the optimal solution for the enlarged dataset. Likewise, when a point k is to be removed from the dataset, its weight is forced to 0, while updating the weights of the remaining points and μ so that the solution obtained with $\alpha_k = 0$ is optimal for the reduced dataset. The online learning follows naturally from the incremental/decremental learning: the new example is added while some old example is removed from the working set.

INCREMENTAL SVM: AN OVERVIEW

Main idea

The basic principle of the incremental SVM [4] is that updates to the state of the example k should keep the remaining examples in their optimal state. In other words, the Kuhn-Tucker (KT) conditions:

$$g_i = -c_i + K_{i,:}x + \mu a_i \begin{cases} \geq 0, & \text{if } x_i = 0 \\ = 0, & \text{if } 0 < x_i < C \\ \leq 0, & \text{if } x_i = C \end{cases} \quad (2)$$

$$\frac{\partial W}{\partial \mu} = a^T x + b = 0 \quad (3)$$

must be maintained for all the examples, except possibly for the current one.

To maintain optimality in practice, one can write out conditions (2)-(3) for the states before and after the update of x_k . By subtracting one from the other the following condition on increments of Δx and Δg is obtained:

$$\begin{bmatrix} \Delta g_k \\ \Delta g_s \\ \Delta g_r \\ 0 \end{bmatrix} = \begin{bmatrix} a_k & K_{ks} \\ a_s & K_{ss} \\ a_r & K_{rs} \\ 0 & a_s^T \end{bmatrix} \underbrace{\begin{bmatrix} \Delta \mu \\ \Delta x_s \end{bmatrix}}_{\Delta r} + \begin{bmatrix} K_{kk}^T \\ K_{ks}^T \\ K_{kr}^T \\ a_k \end{bmatrix} \Delta x_k. \quad (4)$$

The subscript s refer to the examples in the set S of unbounded support vectors, and the subscript r refers to the set R of bounded support vectors (E) and other examples (O). It follows from (2) that $\Delta g_s = 0$. Then lines 2 and 4 of the system (4) can be re-written as:

$$\begin{bmatrix} 0 \\ 0 \end{bmatrix} = \begin{bmatrix} 0 & a_s^T \\ a_s & K_{ss} \end{bmatrix} \Delta s + \begin{bmatrix} a_k \\ K_{ks}^T \end{bmatrix} \Delta x_k. \quad (5)$$

This linear system is easily solved for Δs :

$$\Delta s = \beta \Delta x_k, \quad (6)$$

where

$$\beta = - \underbrace{\begin{bmatrix} 0 & a_s^T \\ a_s & K_{ss} \end{bmatrix}^{-1}}_Q \underbrace{\begin{bmatrix} a_k \\ K_{ks}^T \end{bmatrix}}_\eta \quad (7)$$

is the gradient of the linear manifold of optimal solutions parameterized by x_k .

One can further substitute (6) into the lines 1 and 3 of the system (4) and obtain the following relation:

$$\begin{bmatrix} \Delta g_k \\ \Delta g_r \end{bmatrix} = \gamma \Delta x_k. \quad (8)$$

where

$$\gamma = \begin{bmatrix} \alpha_c & K_{tk} \\ \alpha_r & K_{rr} \end{bmatrix} \beta + \begin{bmatrix} K_{tk} \\ K_{rr} \end{bmatrix} \quad (9)$$

is the gradient of the linear manifold of the gradients of the examples in set R at the optimal solution parameterized by x_k .

Accounting: a systematic account

Notice that all the reasoning in the preceding section is valid only for sufficiently small Δx_k such that the composition of sets S and R does not change. Although computing the optimal Δx_k is not possible in one step, one can compute the largest update Δx_k^{\max} such that composition of sets S and R remains intact. Four cases must be accounted for¹:

1. Some x_i in S reaches a bound (upper or lower one). Let ϵ be a small number. Compute the sets²

$$\begin{aligned} I_+^S &= \{i \in S : \text{sign}(\Delta x_k) \beta_i > \epsilon\} \\ I_-^S &= \{i \in S : \text{sign}(\Delta x_k) \beta_i < -\epsilon\}. \end{aligned}$$

The examples in set I_+^S have positive sensitivity with respect to the current example; that is, their weight would increase by taking a step Δx_k . These examples should be tested for reaching the upper bound C . Likewise, the examples in set I_-^S should be tested for reaching 0. The examples with $-\epsilon < \beta_i < \epsilon$ can be ignored, as they are insensitive to Δx_k . Thus the possible weight updates are:

$$\Delta x_i^{\max} = \begin{cases} C - x_i, & \text{if } i \in I_+^S \\ -x_i, & \text{if } i \in I_-^S, \end{cases}$$

and the largest possible Δx_k^S before one of the elements in S reaches a bound is:

$$\Delta x_k^S = \text{absmin}_{i \in I_+^S \cup I_-^S} \frac{\Delta x_i^{\max}}{\beta_i}, \quad (10)$$

where

$$\text{absmin}(x) := \min_i |x_i| \cdot \text{sign}(x(\arg\min_i |x_i|)).$$

2. Some g_i in R reaches zero. Compute the sets

$$\begin{aligned} I_+^R &= \{i \in R : \text{sign}(\Delta x_k) \gamma_i > \epsilon\} \\ I_-^R &= \{i \in R : \text{sign}(\Delta x_k) \gamma_i < -\epsilon\}. \end{aligned}$$

The examples in set I_+^R have positive sensitivity of the gradient with respect to the weight of the current example. Therefore their (negative)

¹In the original work of Cauwenberghs and Poggio five cases are used but two of them easily fold together.

²Note that $\text{sign}(\Delta x_k)$ is +1 for the incremental and -1 for the decremental cases.

gradients can potentially reach 0. Likewise, gradients of the examples in set \mathcal{I}_-^R are positive but are pushed towards 0 with the changing weight of the current example. Only points in $\mathcal{I}_+^R \cup \mathcal{I}_-^R$ need to be considered for computation of the largest update Δx_k^R :

$$\Delta x_k^R = \text{abmin}_{i \in \mathcal{I}_+^R \cup \mathcal{I}_-^R} \frac{-g_i}{\gamma_i}. \quad (11)$$

3. g_k becomes 0. This case is similar to case 2, except that feasibility test becomes:

$$\text{sign}(\Delta x_k) \gamma_k > \epsilon,$$

and if it holds, the largest update Δx_k^S is computed as:

$$\Delta x_k^S = \frac{-g_k}{\gamma_k}. \quad (12)$$

4. x_k reaches the bound. The largest possible increment is clearly

$$\Delta x_k^S = \begin{cases} C - x_k, & \text{if } x_k \text{ is added} \\ -x_k, & \text{if } x_k \text{ is removed.} \end{cases} \quad (13)$$

Finally, the largest possible update is computed among the four cases:

$$\Delta x_k^{\max} = \text{abmin}((\Delta x_k^S; \Delta x_k^R; \Delta x_k^2; \Delta x_k^1)). \quad (14)$$

The rest of the incremental SVM algorithm essentially consists of repeated computation of the update Δx_k^{\max} , update of the sets S , E and O , update of the state and of the sensitivity parameters β and γ . The iteration stops when either case 3 or case 4 occurs in the increment computation. Computational aspects of the algorithm can be found in [4].

Special case: empty set S

Applying this incremental algorithm leaves open the possibility of an empty set S . This has two main consequences. First, all the blocks with the subscript s vanish from the KT conditions (4). Second, it is impossible to increase the weight of the current example since this would violate the equality constraint of the SVM. As a result, the KT conditions (4) can be written component-wise as

$$\Delta g_k = a_k \Delta \mu \quad (15)$$

$$\Delta g_r = a_r \Delta \mu. \quad (16)$$

One can see that the only free variable is $\Delta \mu$, and $[a_k; a_r]$ plays the role of sensitivity of the gradient with respect to $\Delta \mu$. To select the points from E or O which may enter set S , a feasibility relationship similar to the main case,

can be derived. Resolving (15) for $\Delta\mu$ and substituting the result into (16), we conclude that

$$\Delta g_r = -\frac{a_r}{a_k} \Delta g_k.$$

Then, using the KT conditions (2), the feasible index sets can be defined as

$$I_+ = \{i \in E : -\frac{a_i}{a_k} g_k > \epsilon\} \quad (17)$$

$$I_- = \{i \in O : -\frac{a_i}{a_k} g_k < -\epsilon\} \quad (18)$$

and the largest possible step $\Delta\mu^{\max}$ can be computed as:

$$\Delta\mu^{\max} = \min_{i \in I_+ \cup I_- \cup k} \frac{-g_i}{a_i}. \quad (19)$$

ONLINE SVDD

As it was mentioned in the introduction, the online SVDD algorithm uses the same procedure as the incremental SVM, with the following definitions of the abstract parameters in problem (1): $c = \text{diag}(K)$, $a = y$ and $b = -1$. However, special care needs to be taken of the initialization stage, in order to obtain the initial feasible solution.

Initialization

For the standard support vector classification, an optimal solution for a single point is possible; $x_1 = 0, b = y_1$. In the incremental SVDD the situation is more complicated. The difficulty arises from the fact that the equality constraint $\sum_{i=1}^n a_i x_i = 1$ and the box constraint $0 \leq x_i \leq C$ may be inconsistent; in particular, the constraint cannot be satisfied when fewer than $\lceil \frac{1}{C} \rceil$ examples are available. This initial solution can be obtained by the following procedure:

1. Take the first $\lceil \frac{1}{C} \rceil$ objects, assign them weight C and put them in E .
2. Take the next object k , assign it $x_k = 1 - \lceil \frac{1}{C} \rceil C$ and put it in S .
3. Compute the gradients g_i of all objects, using (2). Compute μ such that for all objects in E the gradient is less than or equal to zero:

$$\mu = -\max_{i \in E} g_i \quad (20)$$

4. Enter the main loop of the incremental algorithm.

+49 30 8825823

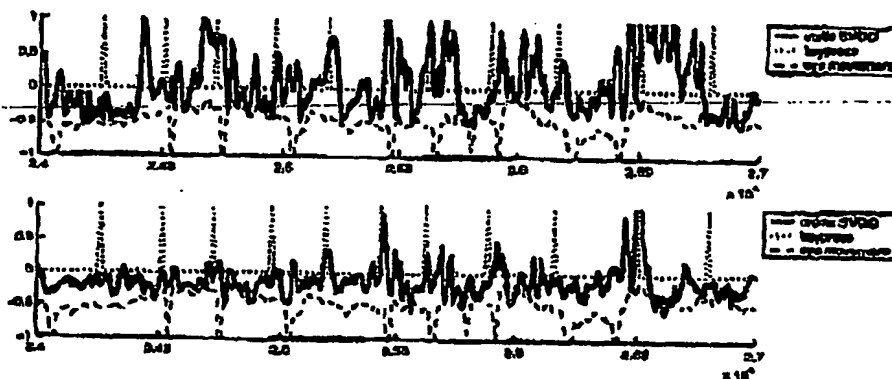


Figure 1: Classification of a time series using a fixed classifier (top) and an online classifier (bottom). The dotted line with the regular peaks are the keystrokes. The noisy solid line indicates the classifier output. The dashed line is the EOG, indicating the activity of the eye (in particular eye-blinks).

Experiments on BCI data

This experiments shows the use of the online novelty detection task on non-stationary time series data. The online SVDD is applied to a BCI (Brain-Computer-Interface) project [2, 3]. A subject was sitting in front of a computer, and was asked to press a key on the keyboard using the left or the right hand. During the experiment, the EEG brain signals of the subject are recorded. From these signals, it is the task to predict which hand will be used for the key press. The first step in the classification task requires a distinction between 'movement' and 'no-movement' which should be made online. The incremental SVDD will be used to characterize the normal activity of the brain, such that special events, like upcoming keystroke movements, are detected.

After preprocessing the EEG signals, at each time point the brain activity is characterized by 21 feature values. The sampling rate was reduced to 10 Hz. A window of 500 time points (thus 5 seconds long) at the start of the time series was used to train an SVDD. In the top plot of figure 1 the output of this SVDD is shown through time. For visualization purposes just a very short, but characteristic part of the time series is shown. The dotted line with the regular single peaks indicates the times at which a key was pressed. The output of the classifier is shown by the solid noisy line. When this line exceeds zero, an outlier, or deviation from the normal situation is detected. The dashed line at the bottom of the graph, shows the muscular activity at the eyes. The large spikes indicate eye blinks, which are also detected as outliers. It appears that the output of the static classifier through time is very noisy. Although it detects some of the movements and eye blinks, it also generates many false alarms.

In the bottom plot of figure 1 the output of the online SVDD classifier is

+49 30 8825823

TABLE 1: TEST CLASSIFICATION ERRORS ON THE USPS DATASET, USING A SUPPORT VECTOR CLASSIFIER (RBF KERNEL, $\sigma^2 = 0.3 \cdot 256$) WITH JUST M OBJECTS.

M	50	100	150	200	250	300	350	∞
error (%)	25.41	6.88	4.68	4.48	4.43	4.38	4.29	4.25

shown. Here again, an output above zero indicates that an outlier is detected. It is clear that the online-version generates less false alarms, because it follows the changing data distribution. Although the detection is far from perfect, as can be observed, many of the keystrokes are indeed clearly detected as outliers. It is also clear that the method is easily triggered by the eye blinks. Unfortunately the signal is very noisy, and it is hard to quantify the exact performance for these methods on this data.

ONLINE LEARNING IN LARGE DATASETS

To make the SVM learning applicable to very large datasets, the classifier has to be constrained to have a limited number of objects in memory. This is, in principle, exactly what an online classifier with fixed window size M does. The only difference is that removing the oldest object is not useful in this application because the same result is achieved as if the learning had been done on the last M objects. Instead, the "least relevant" object needs to be removed during each window advancement. A reasonable criterion for relevance seems to be the value of the weight. In the experiment presented below the example with the smallest weight is removed from the working set.

Experiments on the USPS data

The dataset is the standard US Postal Service dataset, containing 7291 training and 2007 images of handwritten digits, size 16×16 [19]. On this 10 class dataset 10 support vector classifiers with a RBF kernel, $\sigma^2 = 0.3 \cdot 256$ and $C = 100$, were trained³. During the evaluation of a new object, it is assigned to the class corresponding to the classifier with the largest output. The total classification error on the test set for different window sizes M is shown in table 1.

One can see that the classification accuracy deteriorates marginally (by about 10%) until the working size of 150, which is about 2% of the data. Clearly, by discarding "irrelevant" examples, one removes potential support vectors that cannot be recovered at a later stage. Therefore it is expected that performance of the limited memory classifier would be worse than that of an unrestricted classifier. It is also obvious that no more points than the number of support vectors are eventually needed, although the latter number is not known in advance. The average number of support vectors per each unrestricted 2-class classifier in this experiment is 274. Therefore the results above can be interpreted as reducing the storage requirement by 46% from

³The best model parameters as reported in [19] were used.

the minimal at the cost of 10% increase of classification problem.

Notice that the proposed strategy differs from the caching strategy, typical for many SVM^{light}-like algorithms [6, 8, 5], in which kernel products are re-computed if the examples are found missing in the fixed-size cache and the accuracy of the classifier is not sacrificed. Our approach constitutes a trade-off between accuracy and computational load because kernel products never need to be re-computed. It should be noted, however, that computational cost of re-computing the kernels can be very significant, especially for the problems with complicated kernels such as string matching or convolution kernels.

CONCLUSIONS

Based on revised version of the incremental SVM, we have proposed: (a) an online SVDD algorithm which, unlike all previous extensions of incremental SVM, deals with an unsupervised learning problem, and (b) a fixed-memory training algorithm for the classification SVM which allows to limit the memory requirement for storage of the kernel matrix at the expense of classification performance. Experiments on novelty detection in non-stationary time series and on the USPS dataset demonstrate feasibility of both approaches. More detailed comparisons with other subsampling techniques for limited-memory learning will be carried out in future work.

Acknowledgements

This research was partially supported through a European Community Marie Curie Fellowship and BMBF FKZ 01IBB02A. We would like to thank K.-R. Müller and B. Blankertz for fruitful discussions and the use of BCI data. The authors are solely responsible for information communicated and the European Commission is not responsible for any views or results expressed.

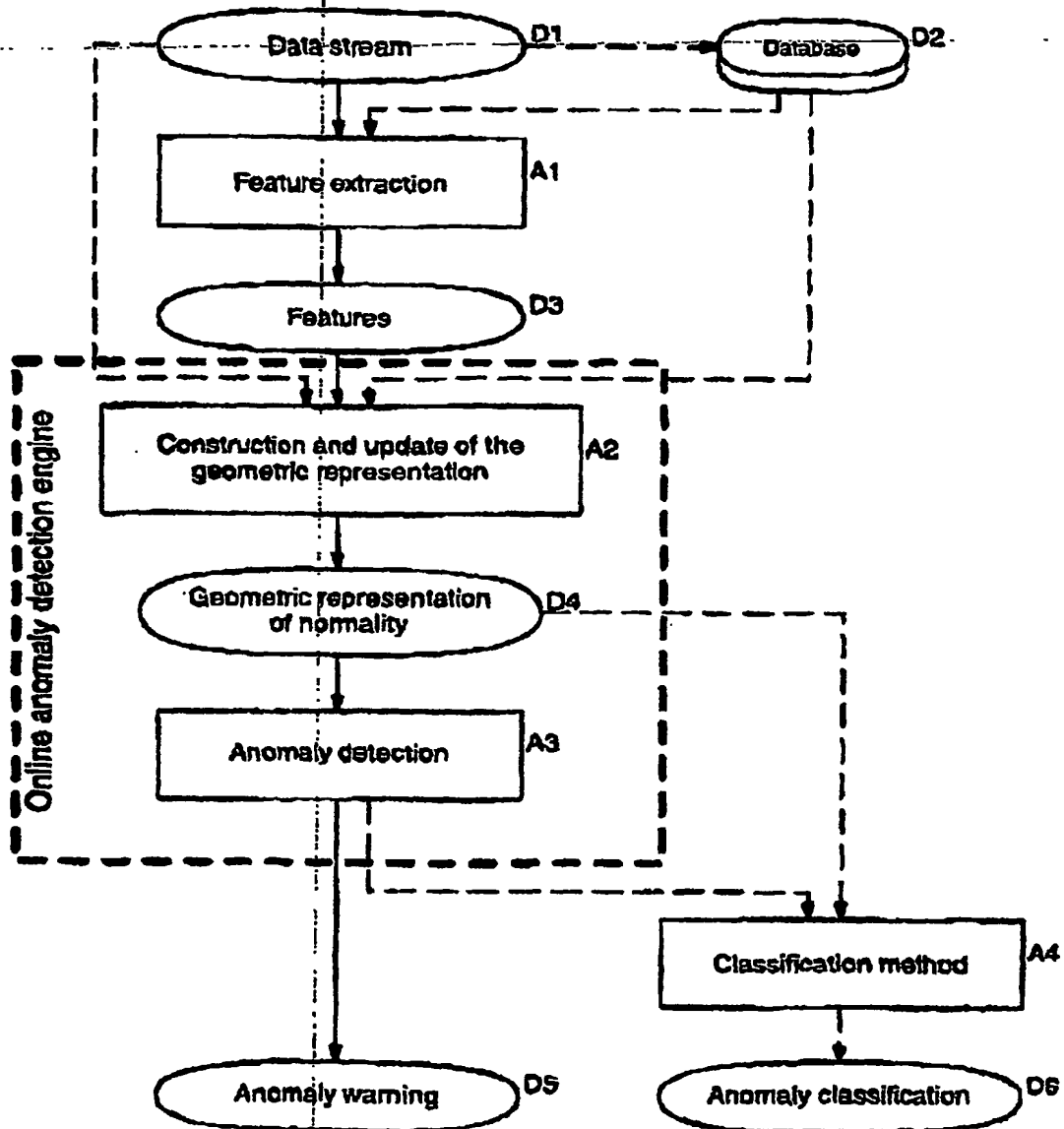
REFERENCES

- [1] D. Achlioptas, F. McSherry and B. Schölkopf, "Sampling Techniques for Kernel Methods," in T. Diettrich, S. Becker and Z. Ghahramani (eds.), *Advances in Neural Information Processing Systems*, 2002, vol. 14, pp. 335-341.
- [2] B. Blankertz, G. Curio and K.-R. Müller, "Classifying Single Trial EEG: Towards Brain Computer Interfacing," in T. G. Diettrich, S. Becker and Z. Ghahramani (eds.), *Advances in Neural Inf. Proc. Systems (NIPS 01)*, 2002, vol. 14, pp. 157-164.
- [3] B. Blankertz, G. Dornhege, O. Schärer, R. Krepki, J. Kohlmorgen, K.-R. Müller, V. Kunzmann, F. Loech and G. Curio, "BCI bit rates and error de-

- tection for fast-pace motor commands based on single-trial EEG analysis," *IEEE Transactions on Rehabilitation Engineering*, 2003, accepted.
- [4] G. Cauwenberghs and T. Poggio, "Incremental and decremental support vector machine learning," in *Neural Information Processing Systems*, 2000.
- [5] R. Collobert and S. Bengio, "SVMTool: Support vector machines for large-scale regression problems," *Journal of Machine Learning Research*, vol. 1, pp. 143-160, 2001.
- [6] T. Joachims, "Making Large-Scale SVM Learning Practical," in B. Schölkopf, C. Burges and A. Smola (eds.), *Advances in Kernel Methods - Support Vector Learning*, Cambridge, MA: MIT Press, 1999, pp. 169-184.
- [7] J. Klivinen, A. Smola and R. Williamson, "Online learning with kernels," in T. G. Diettrich, S. Becker and Z. Ghahramani (eds.), *Advances in Neural Inf. Proc. Systems (NIPS 01)*, 2001, pp. 785-792.
- [8] P. Laskov, "Feasible direction decomposition algorithms for training support vector machines," *Machine Learning*, vol. 46, pp. 315-349, 2002.
- [9] J. Ma, J. Thaler and S. Perkins, "Accurate online support vector regression," <http://nie-www.lanl.gov/~jt/Papers/soavr.pdf>.
- [10] M. Martin, "On-line Support Vector Machines for function approximation," Techn. report, Universitat Politècnica de Catalunya, Departament de Llenguatges i Sistemes Informàtics, 2002.
- [11] M. Moya and D. Husb, "Network constraints and multi-objective optimization for one-class classification," *Neural Networks*, vol. 9, no. 3, pp. 463-474, 1996.
- [12] L. Ralaivola and F. d'Alché Buc, "Incremental Support Vector Machine Learning: A Local Approach," *Lecture Notes in Computer Science*, vol. 2130, pp. 322-329, 2001.
- [13] S. Rüping, "Incremental learning with support vector machines," Techn. Report TR-18, Universität Dortmund, SFB476, 2002.
- [14] B. Schölkopf, J. Platt, J. Shawe-Taylor, A. Smola and R. Williamson, "Estimating the support of a high-dimensional distribution," *Neural Computation*, vol. 13, no. 7, pp. 1443-1471, 2001.
- [15] B. Schölkopf, A. Smola, R. Williamson and P. Bartlett, "New Support Vector Algorithms," *Neural Computation*, vol. 12, pp. 1207 - 1245, 2000, also NeuroCOLT Technical Report NC-TR-1998-031.
- [16] A. Smola and B. Schölkopf, "Sparse greedy matrix approximation for machine learning," in P. Langley (ed.), *Proc. ICML'00*, San Francisco: Morgan Kaufmann, 2000, pp. 911-918.
- [17] N. A. Syed, H. Liu and K. K. Sung, "Incremental learning with support vector machines," in *SVM workshop, IJCAI*, 1999.
- [18] D. Tax and R. Duin, "Uniform object generation for optimizing one-class classifiers," *Journal for Machine Learning Research*, pp. 155-173, 2001.
- [19] V. Vapnik, *Statistical Learning Theory*, New York: Wiley, 1998.

EPO-BERLIN

19-08-2003

**Fig. 1: SCHEMATIC DIAGRAM OF THE SYSTEM.**

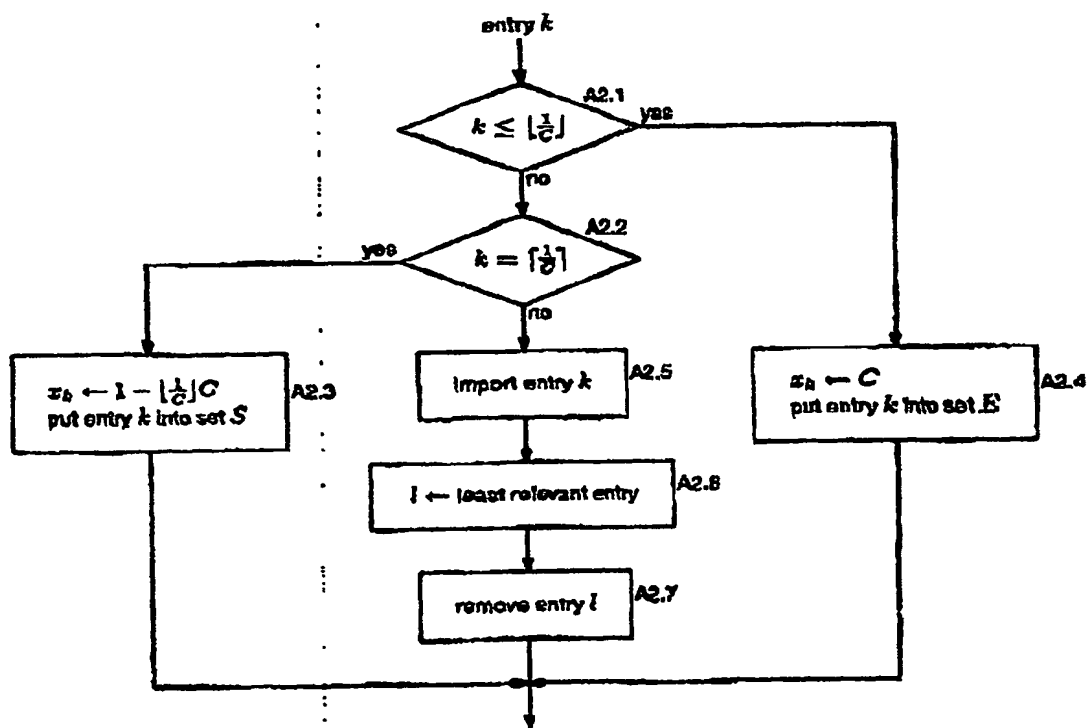


Fig. 2: CONSTRUCTION AND UPDATE OF THE GEOMETRIC REPRESENTATION.

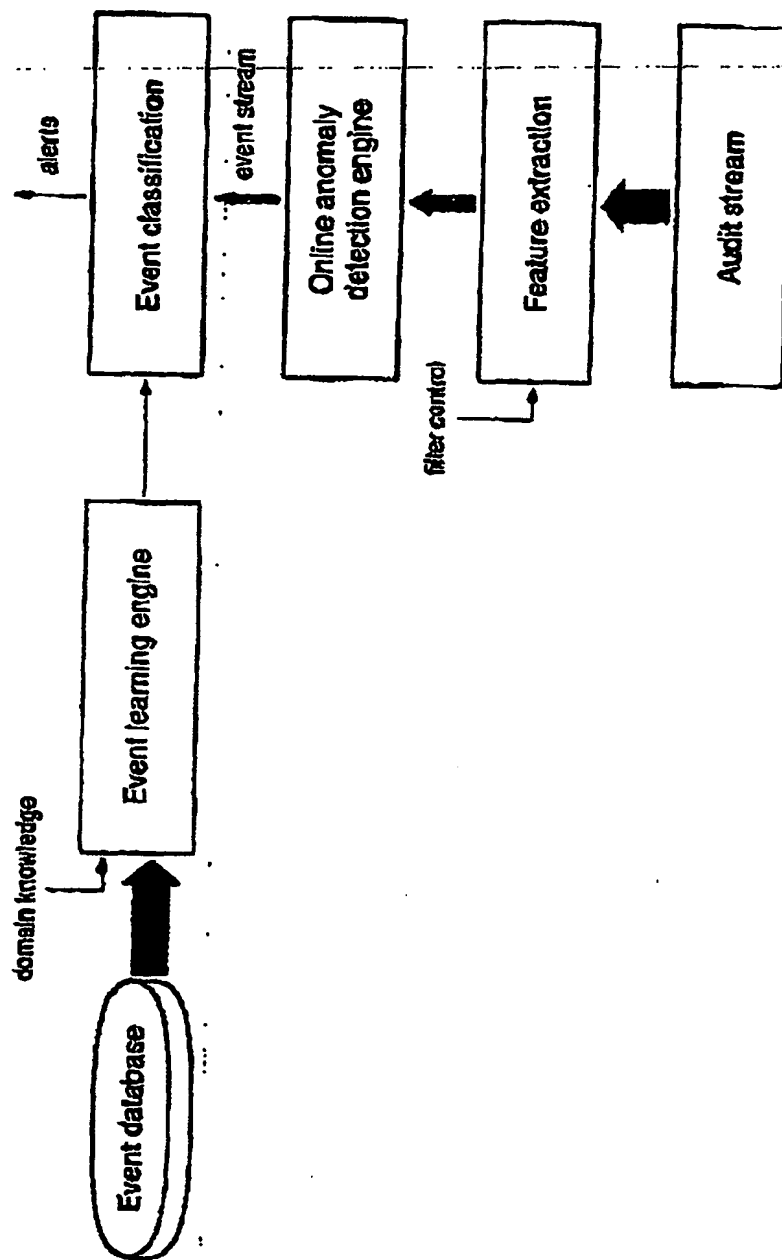


Fig. 3: SYSTEM FOR DETECTION AND CLASSIFICATION OF COMPUTER INTRUSIONS.

THIS PAGE BLANK (USPTO)